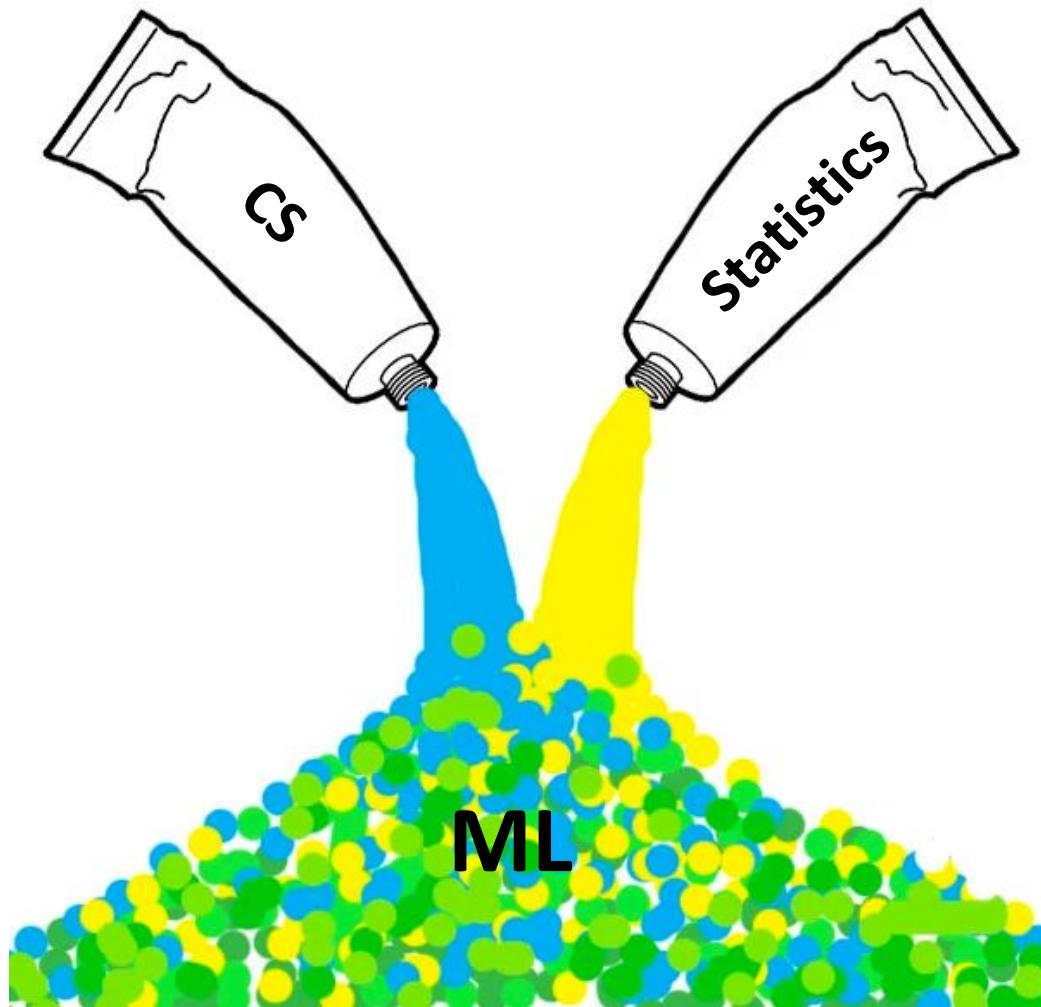




# Machine learning? Statistics?



Jeanette Mumford  
University of Wisconsin - Madison  
Center for Healthy Minds



# Main ideas

- Is machine learning just fancy statistics?
- Is machine learning better?
- How are people using ML in neuroscience?
- Things to consider when transitioning from voxelwise stats to ML

# Main ideas

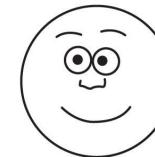
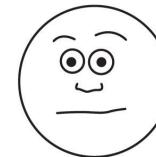
- Is machine learning just fancy statistics?
- Is machine learning better?
- How are people using ML in neuroscience?
- Things to consider when transitioning from voxelwise stats to ML

# Simplified comparison

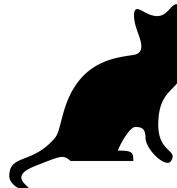
- Statistics
  - Inference
  - Mechanism
- Machine learning
  - Prediction

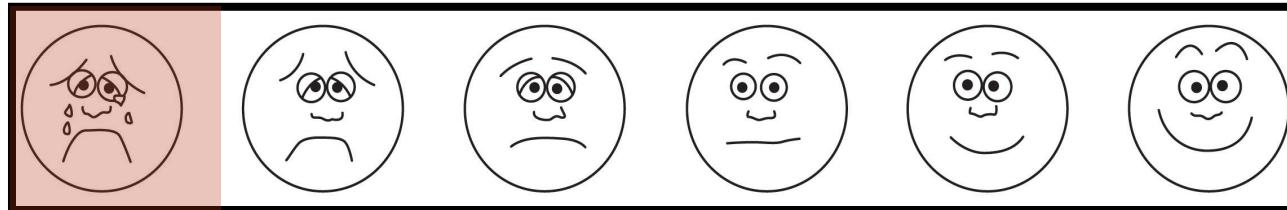
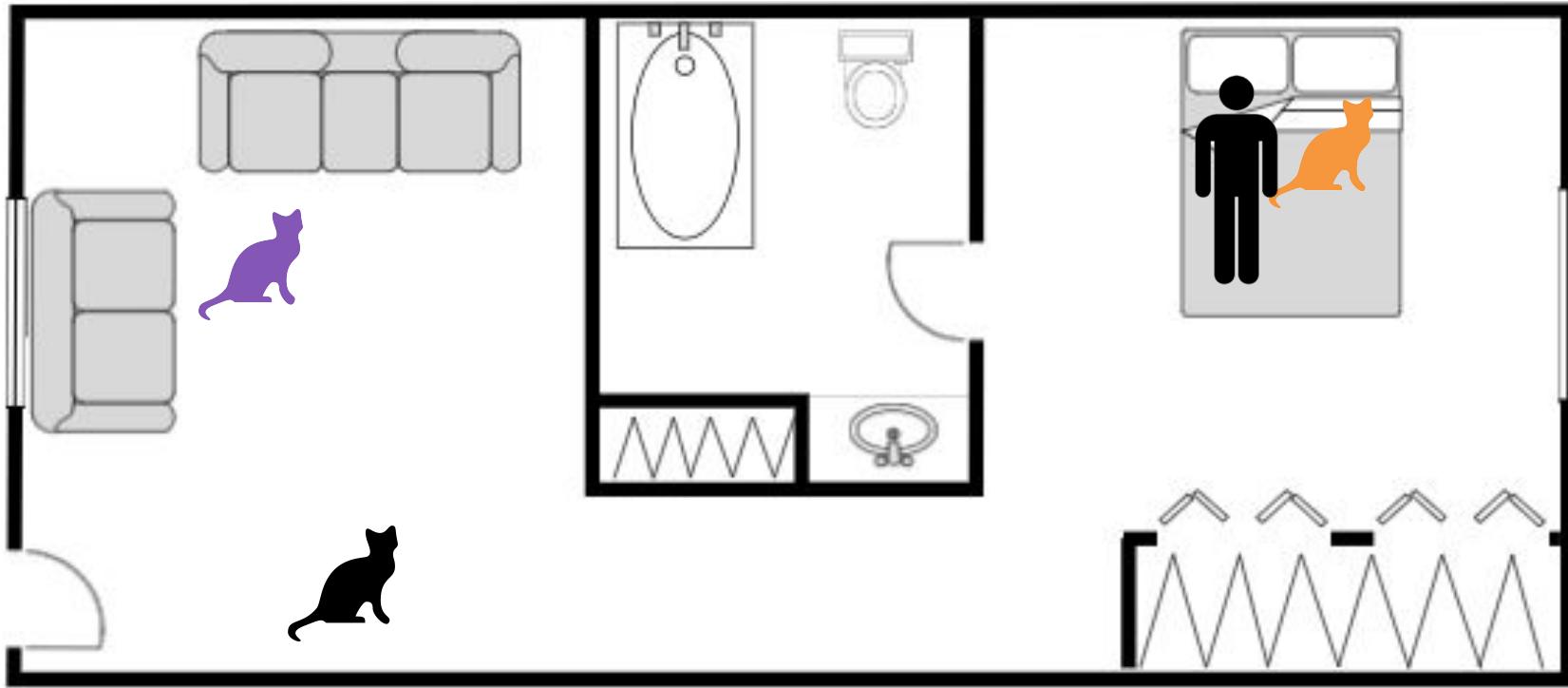
# Example

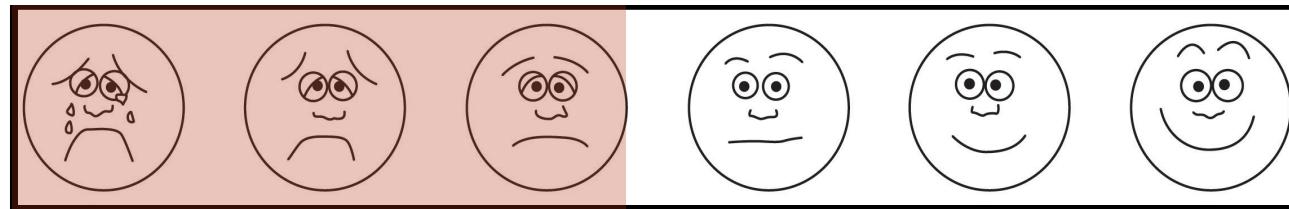
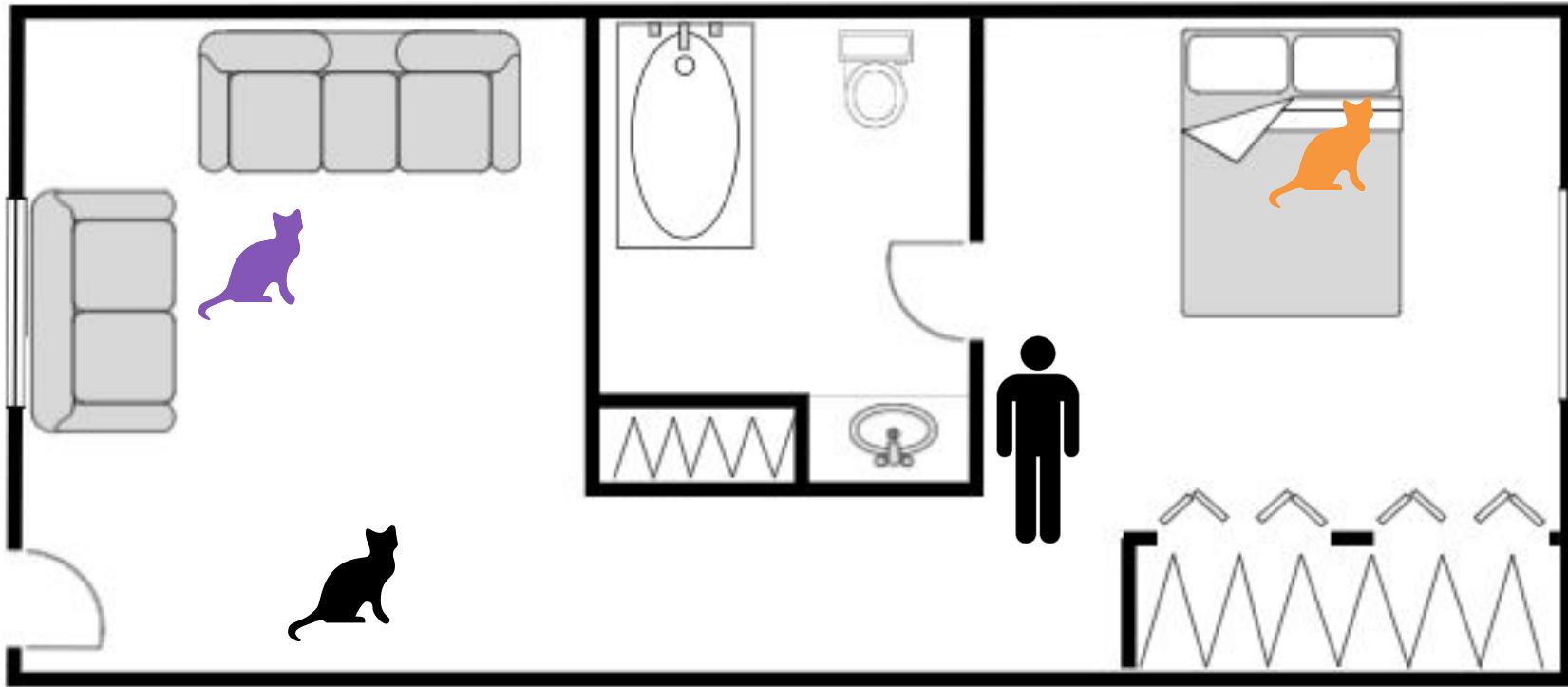
- Measure 1: My current level of happiness over the course of the day

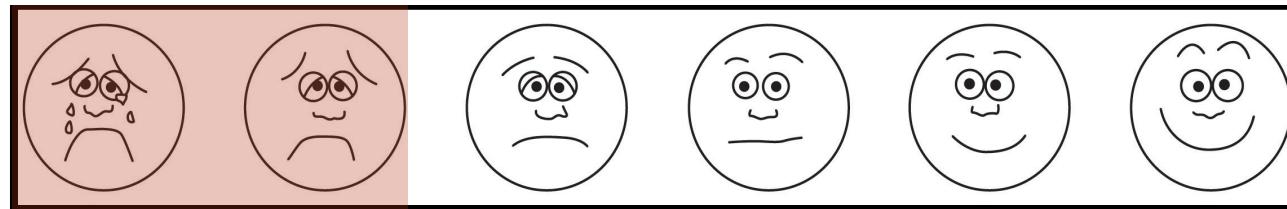
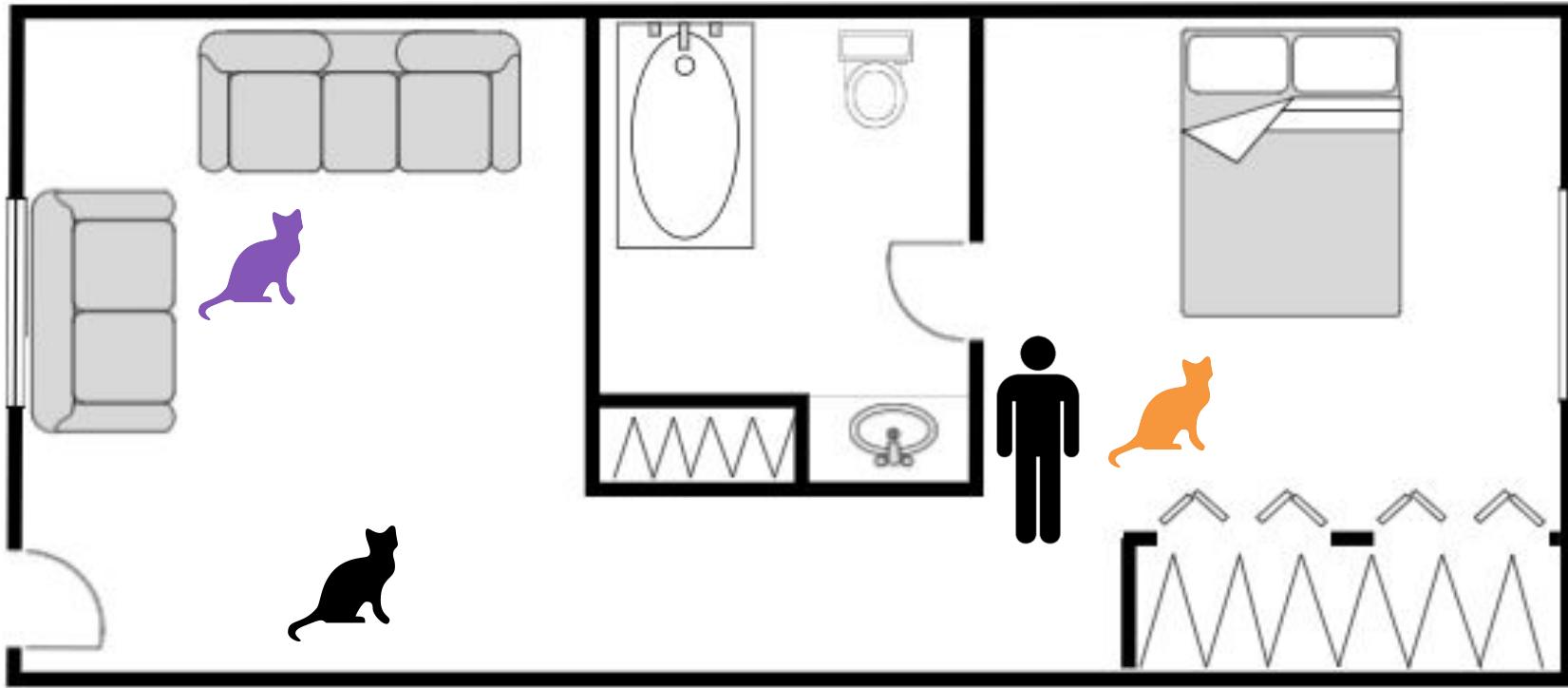


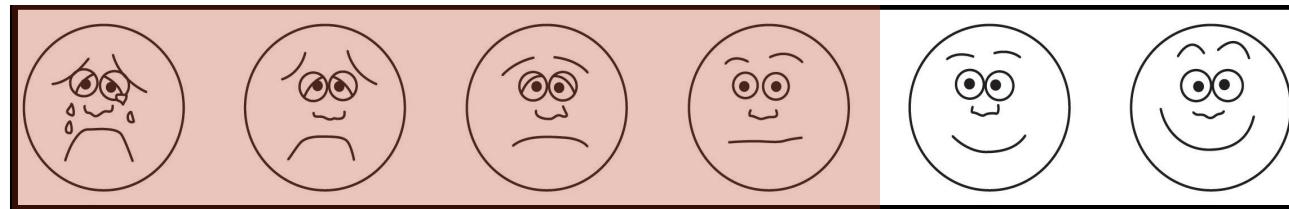
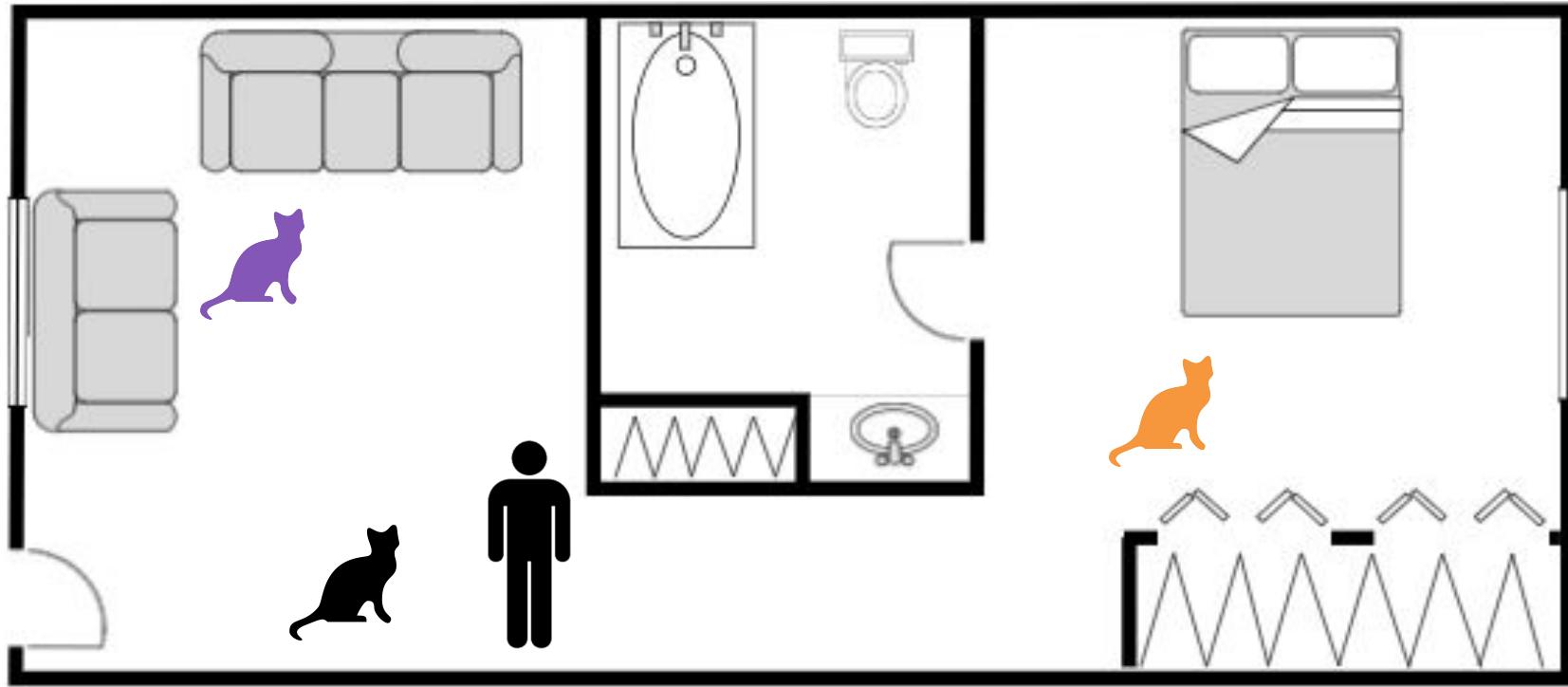
- Measures 2-4: My proximity to each of three different cats

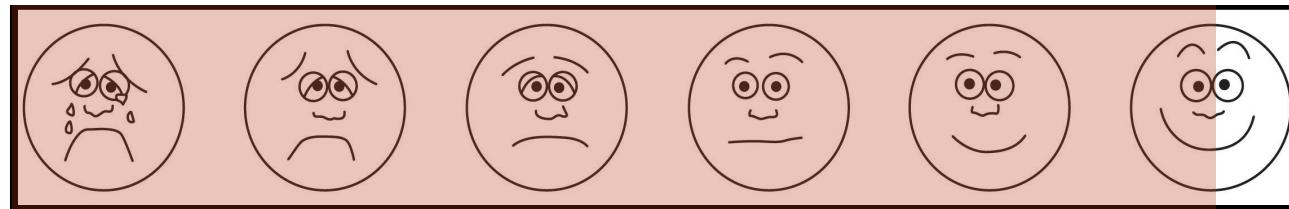
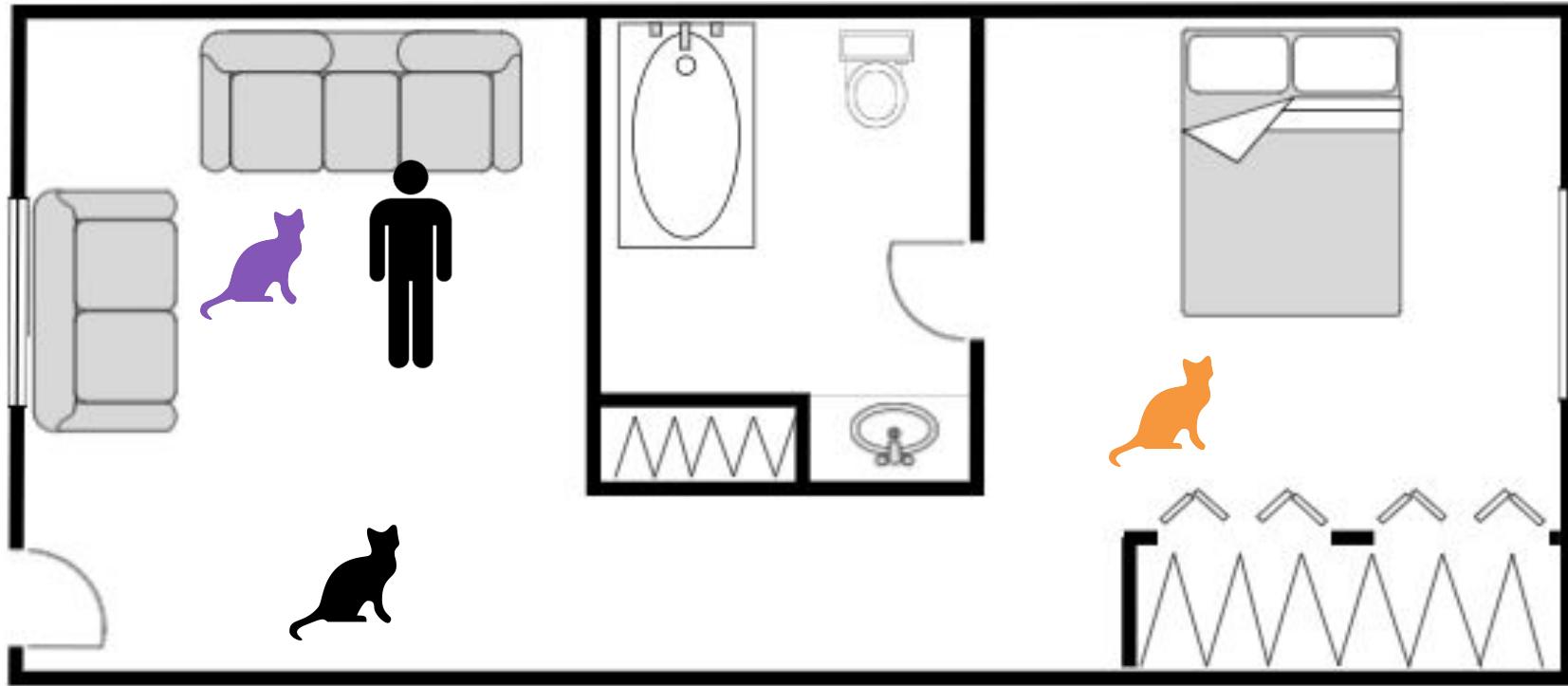












# Statistics questions

- What is the relationship between each cat and happiness?
- Are those relationships statistically significant?

Call:

```
lm(formula = happiness ~  +  + , data = dat1)
```

Residuals:

Min	1Q	Median	3Q	Max
-5.8719	-1.4738	0.1606	1.3308	5.7411

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	5.60892	0.51750	10.839	< 2e-16 ***
	-0.16764	0.03884	-4.316	2.17e-05 ***
	-0.01025	0.06824	-0.150	0.88067
	0.17290	0.06410	2.697	0.00739 **
---				

Signif. codes: 0 ‘\*\*\*’ 0.001 ‘\*\*’ 0.01 ‘\*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 1.973 on 296 degrees of freedom

Multiple R-squared: 0.1071, Adjusted R-squared: 0.09803

F-statistic: 11.83 on 3 and 296 DF, p-value: 2.431e-07

# Machine learning

- If I was to collect a new data set, could the model on the previous slide be used to predict my happiness?
  - Do the betas generalize
  - Just part of a 2-fold cross validation

Mean absolute error = 1.7

Root of mean square error = 2.15

Actual Happiness

12

8

4

0

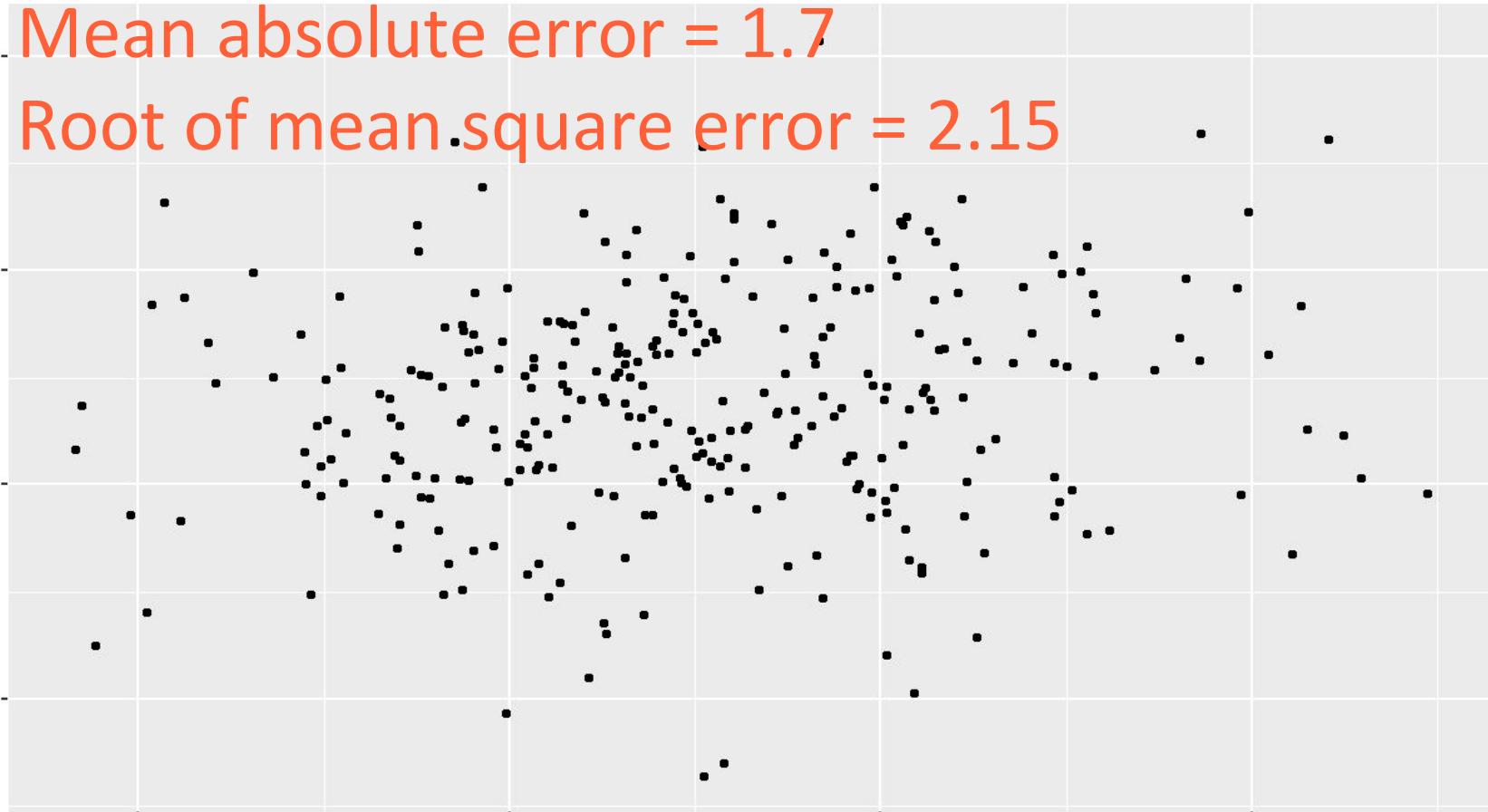
4

5

6

7

Predicted Happiness



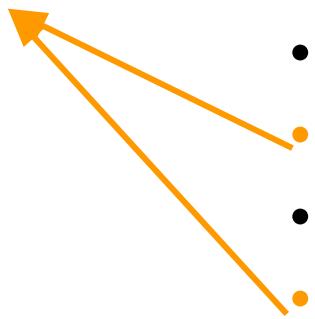
# Deeper comparison

- Machine learning
  - User friendly
  - Good predictions
  - Automated
  - Black box
  - Overfitting worry
  - Generalizability
  - Lots of data
- Statistical Modeling
  - Parsimony
  - Interpretability
  - Knowledge creation
  - Inference
  - Population
  - Generalizability
  - Little data

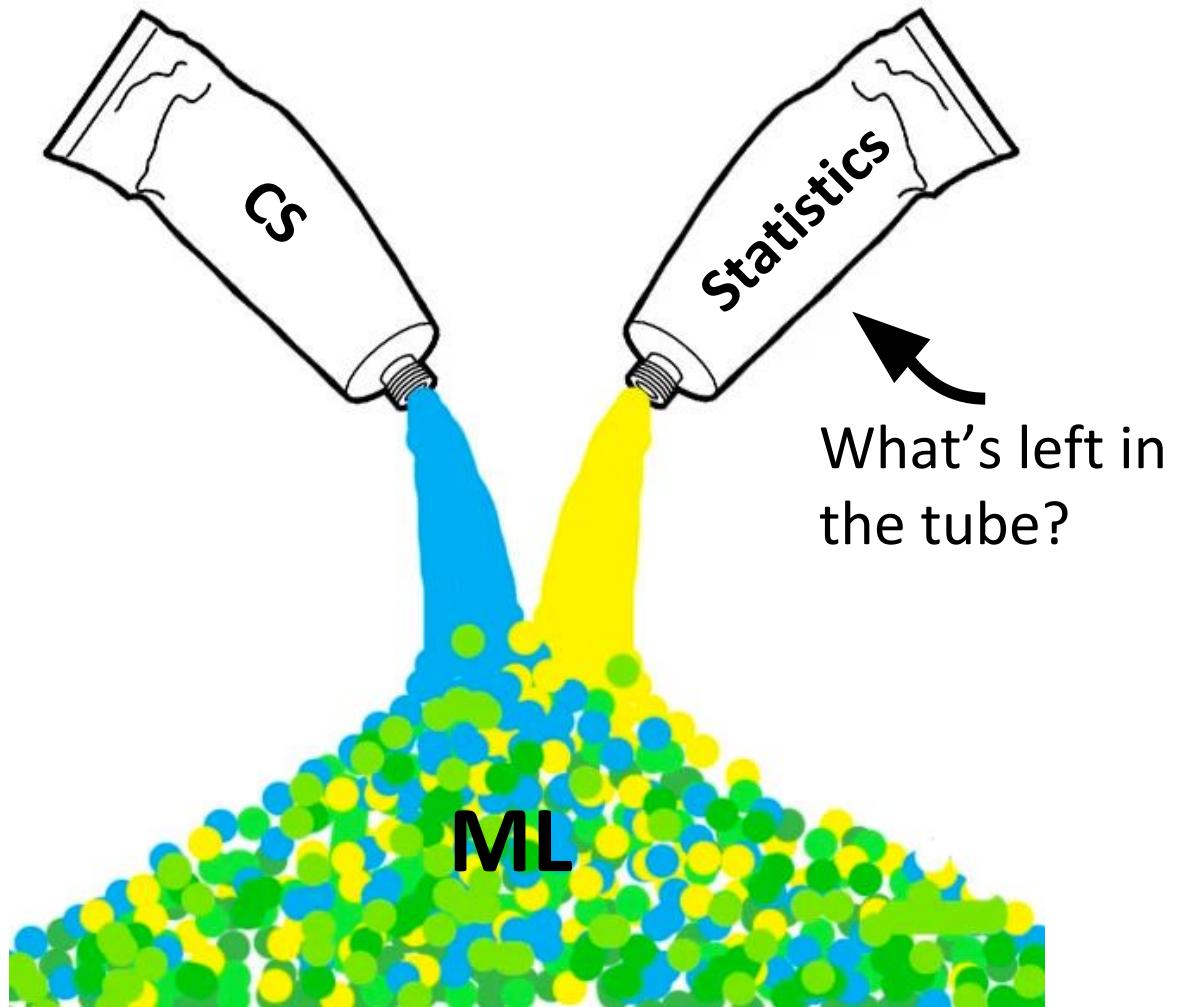
From Brian Caffo's Machine Learning vs Classical Statistics video  
<https://www.youtube.com/watch?v=U0XIBBuJaI4&t=101s>

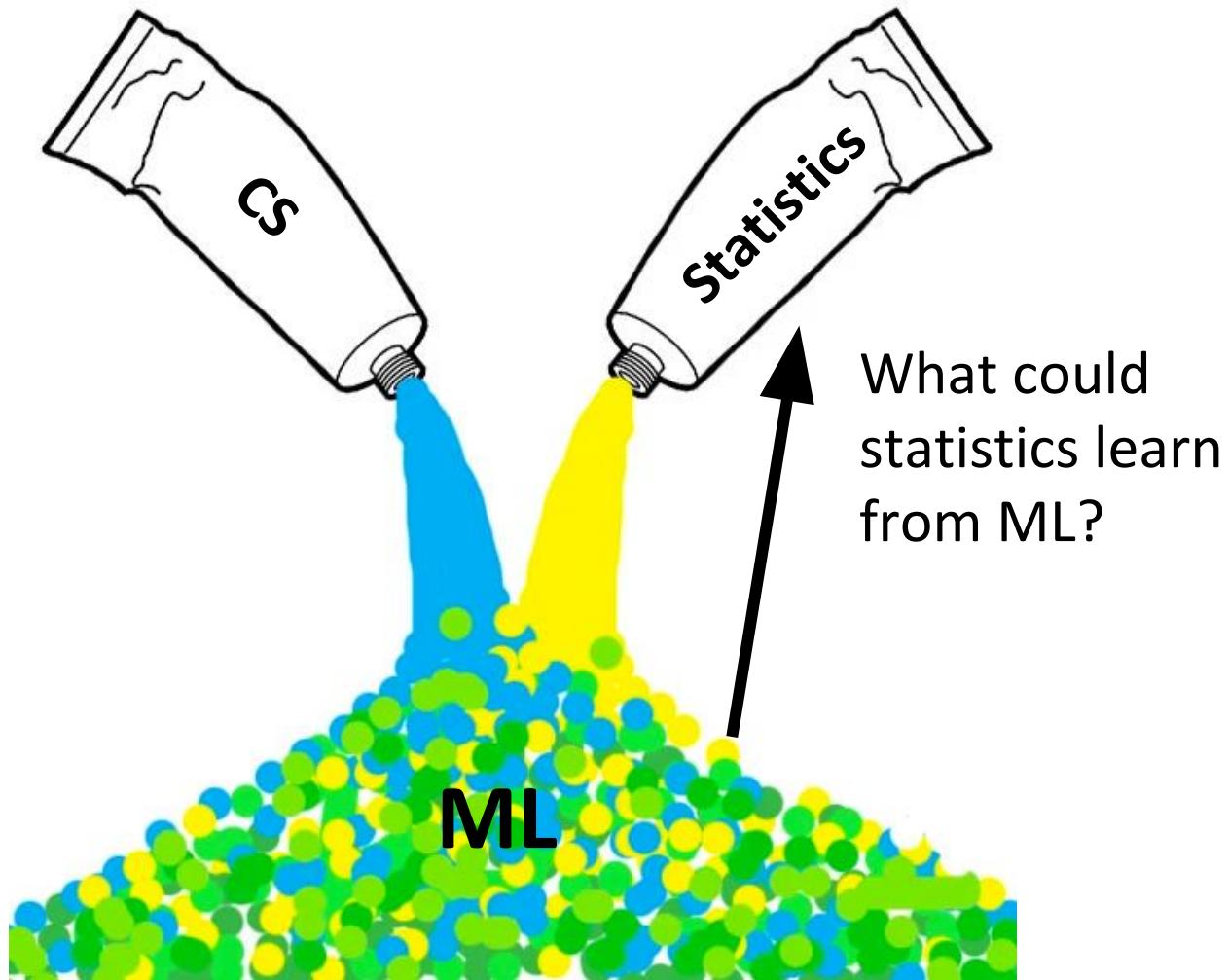
# Deeper comparison

- Machine learning
  - User friendly
  - Good predictions
  - Automated
  - Black box
  - Overfitting worry
  - Generalizability
  - Lots of data
- Statistical Modeling
  - Parsimony
  - Interpretability
  - Knowledge creation
  - Inference
  - Population
  - Generalizability
  - Little data



From Brian Caffo's Machine Learning vs Classical Statistics video  
<https://www.youtube.com/watch?v=U0XIBBuJaI4&t=101s>



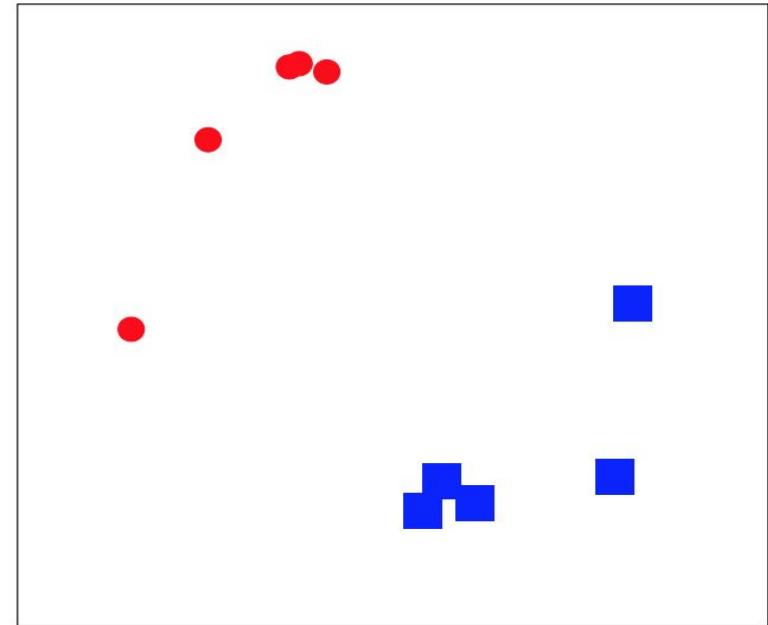


# Comparison

- ML
  - Papers driven by conferences
    - Faster paced
    - Creates a synergy in the field
  - Has dominated some stats topics that are currently ignored in stats
- Statistics
  - Contains important topics that could be useful in ML
    - Worry about uncertainty

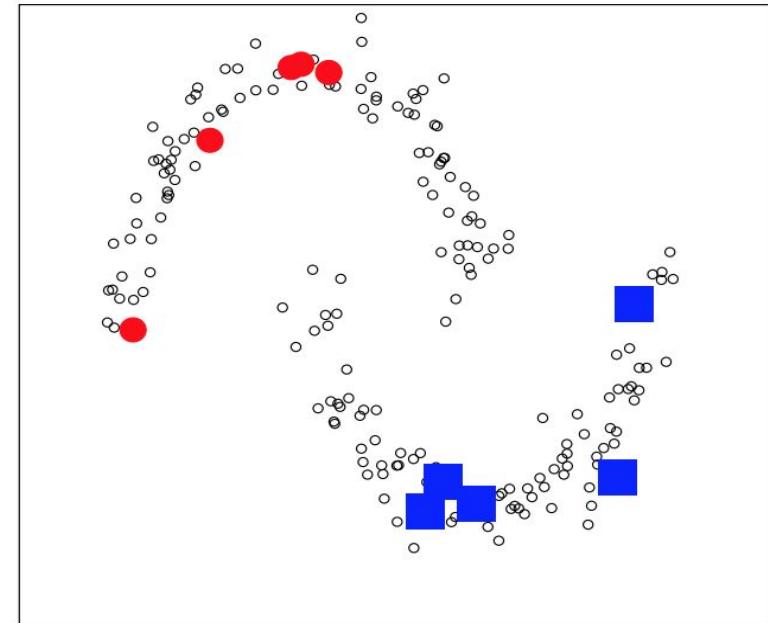
# Example: Semisupervised inference

- You could try to use these data to make predictions



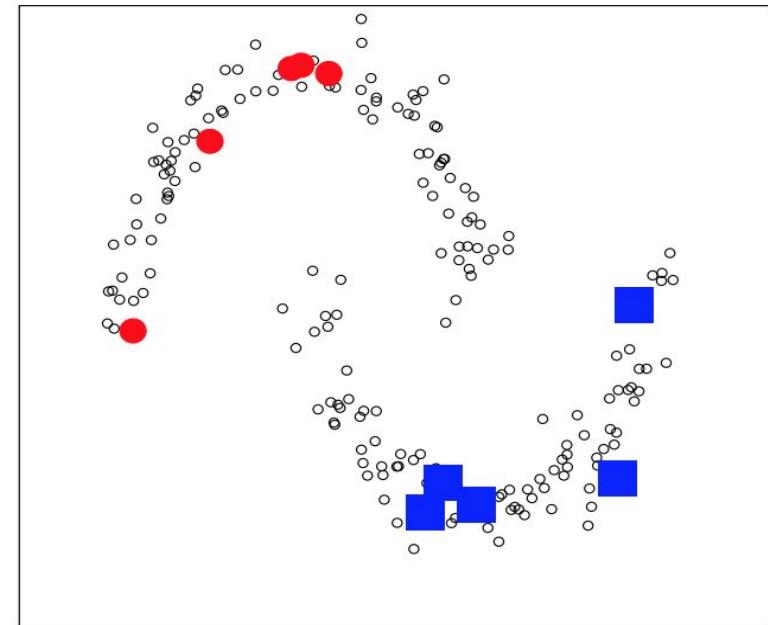
# Example: Semisupervised inference

- You could try to use these data to make predictions
- Unlabeled data may help!



# Example: Semisupervised inference

- You could try to use these data to make predictions
- Unlabeled data may help!
- Inference helps evaluate the assumption
  - Should unlabeled data be used?



More useful to think about how the two fields can benefit from each other



# Although....

- You must consider these differences if you're making career path choices

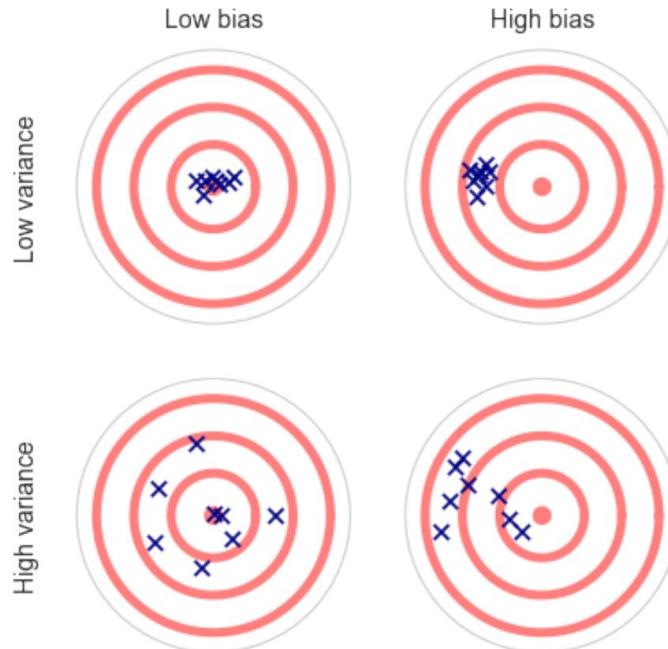
# Main ideas

- Is machine learning just fancy statistics?
- **Is machine learning better?**
- How are people using ML in neuroscience?
- Things to consider when transitioning from voxelwise stats to ML

# Stronger conclusions

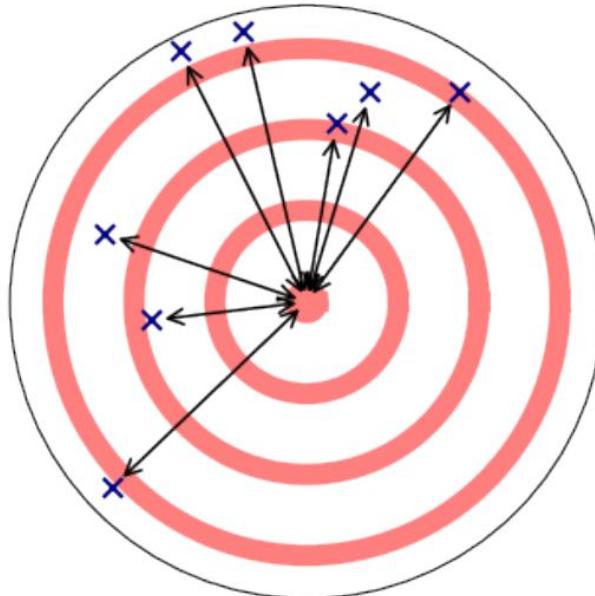
- Asks deeper questions about model coefficients
  - $Y = B_0 + B_1 X_1 + B_2 X_2$ 
    - Regression focus: The model
    - ML focus: Do the betas generalize?
- Get closer to understanding the true usefulness of our theory

# Bias-variance tradeoff

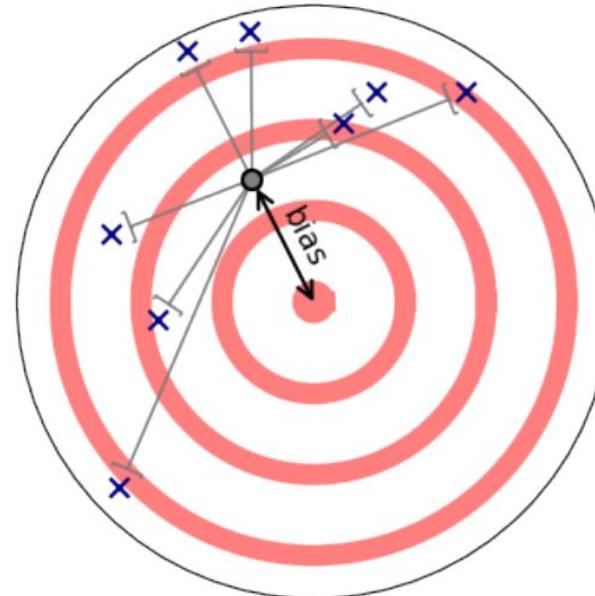


# Bias-variance tradeoff

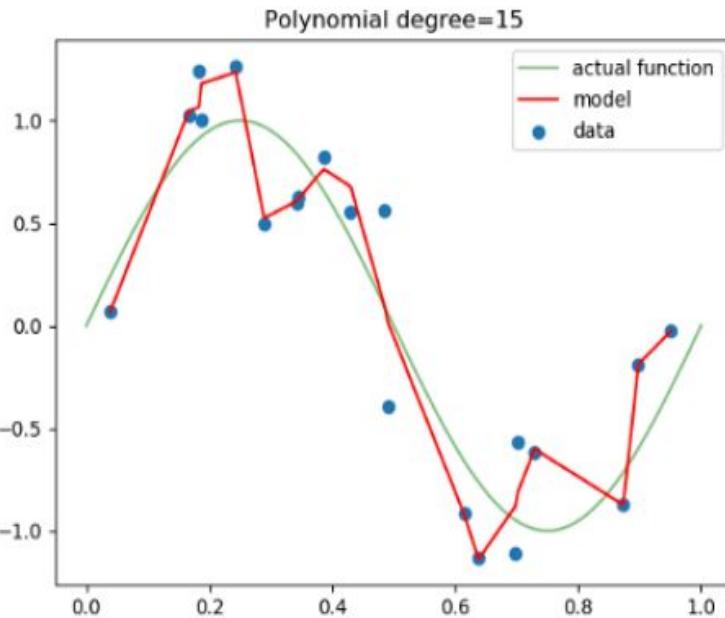
Sum of squared errors



Bias-variance decomposition

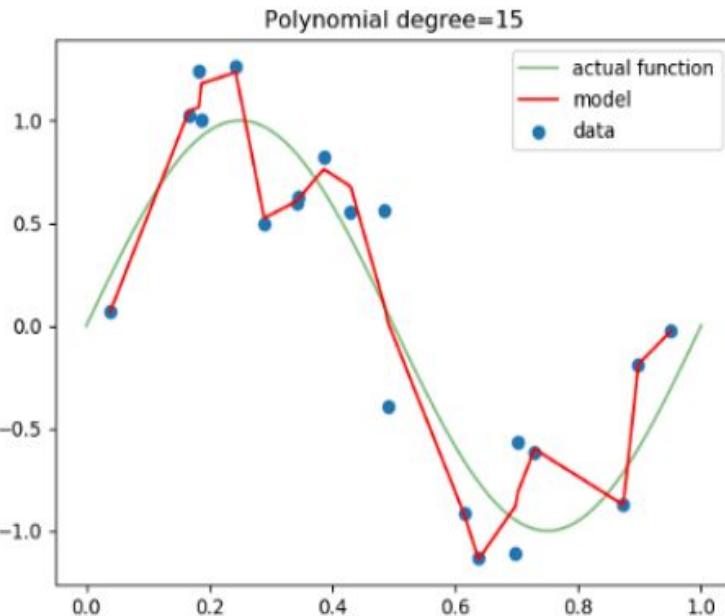


# What does this look like?

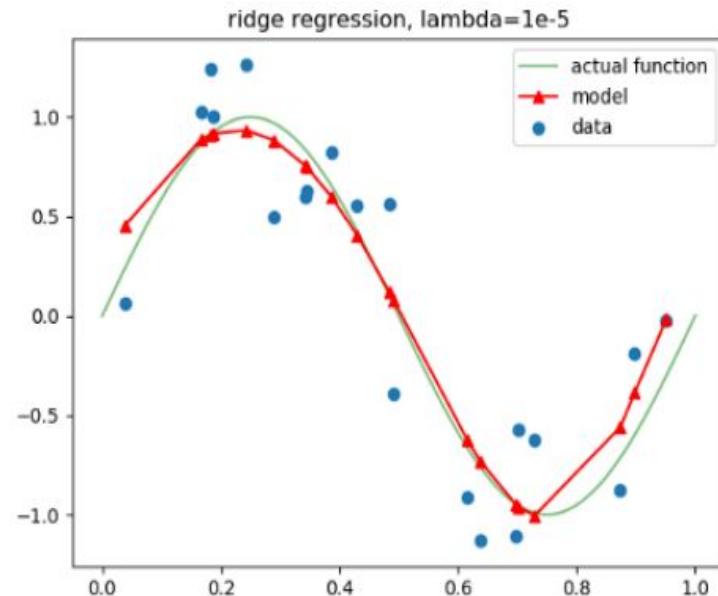


Low bias and high variance regression fit  
(no regularization)

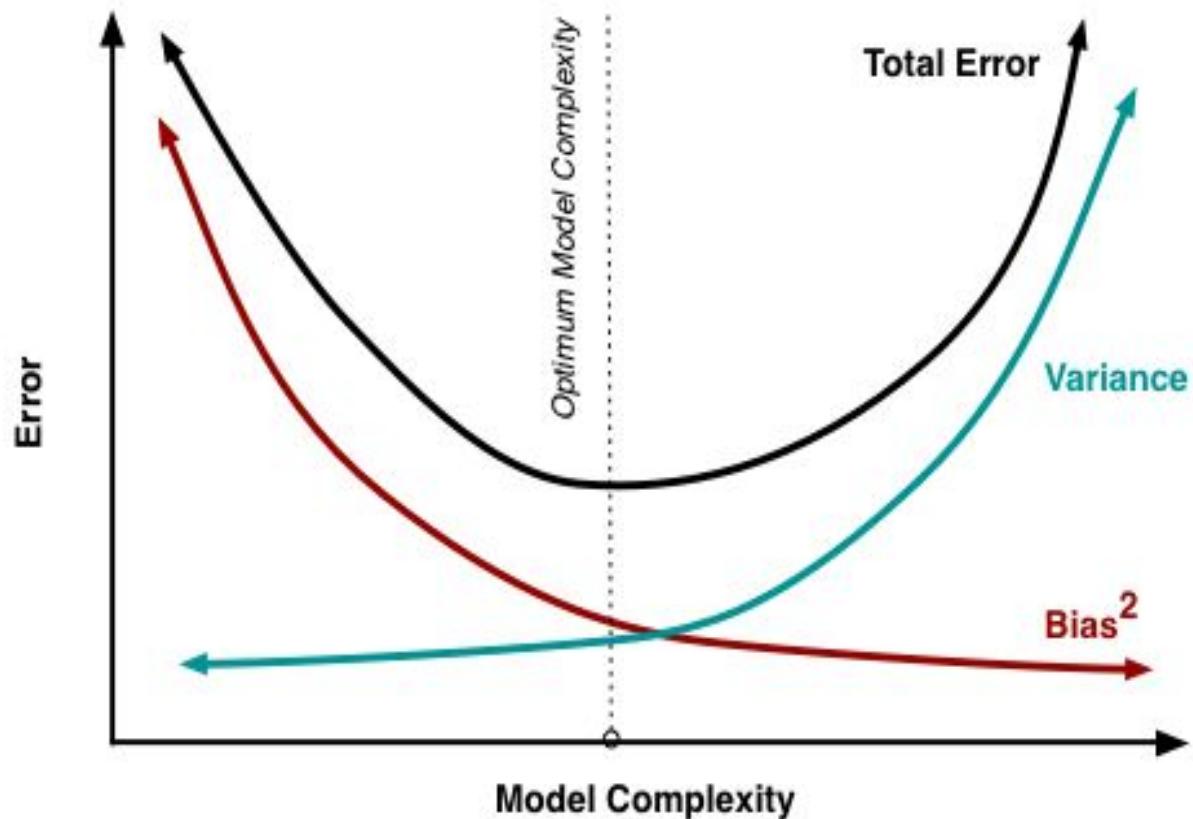
# What does this look like?

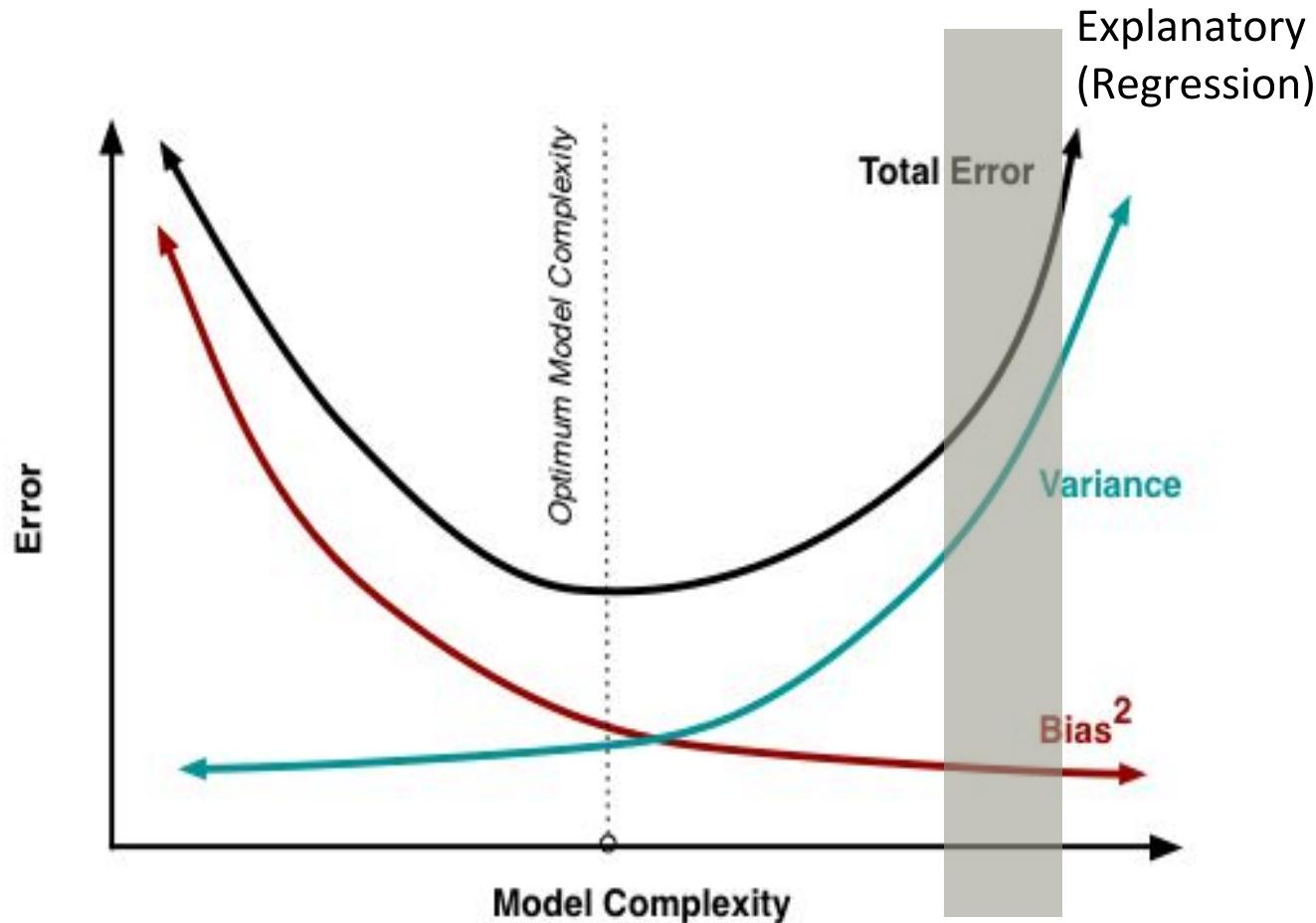


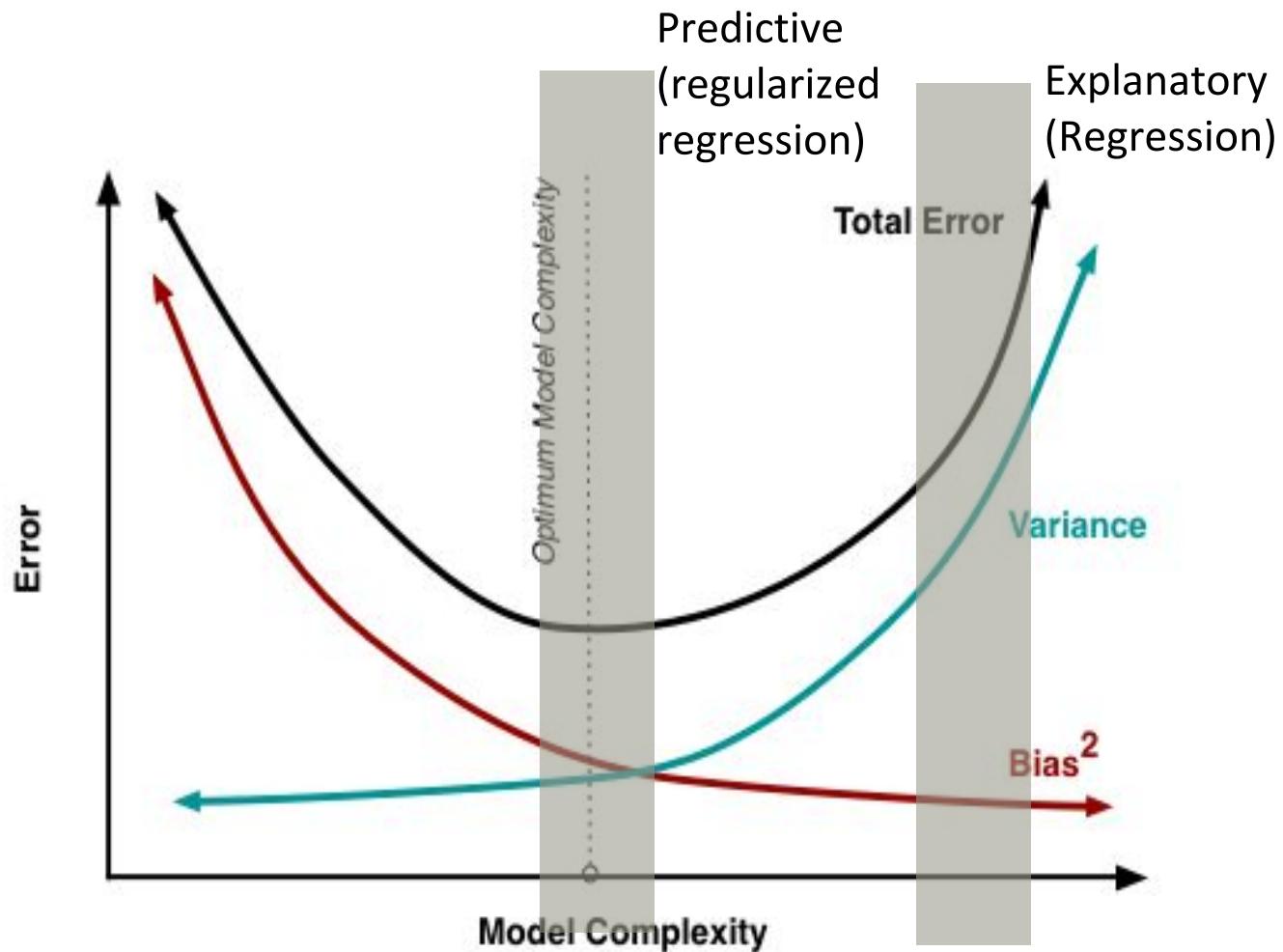
Low bias and high variance regression fit  
(no regularization)



Bias some betas to 0 (ridge regression) to lower the variance -> better predictions







# Voxelwise/ROI analyses aren't reliable

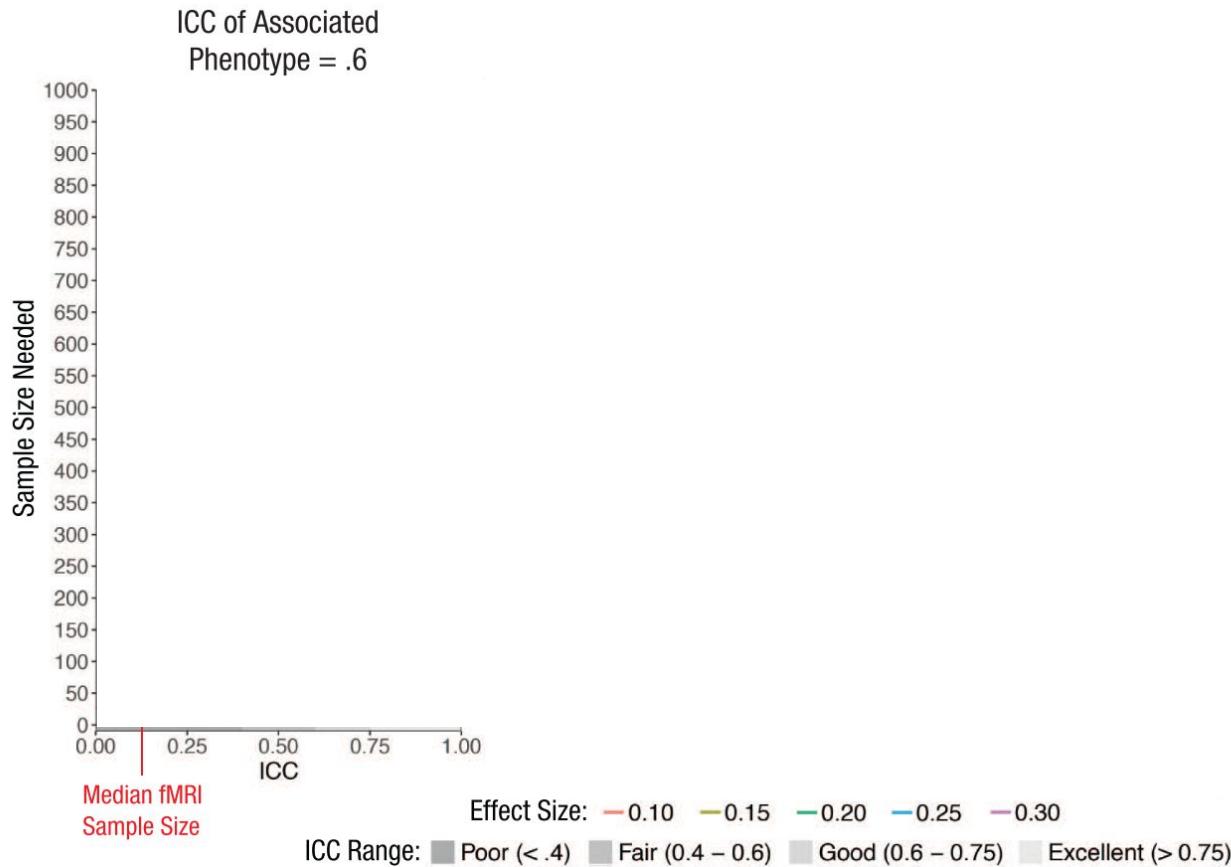
Research Article

## What Is the Test-Retest Reliability of Common Task-Functional MRI Measures? New Empirical Evidence and a Meta-Analysis

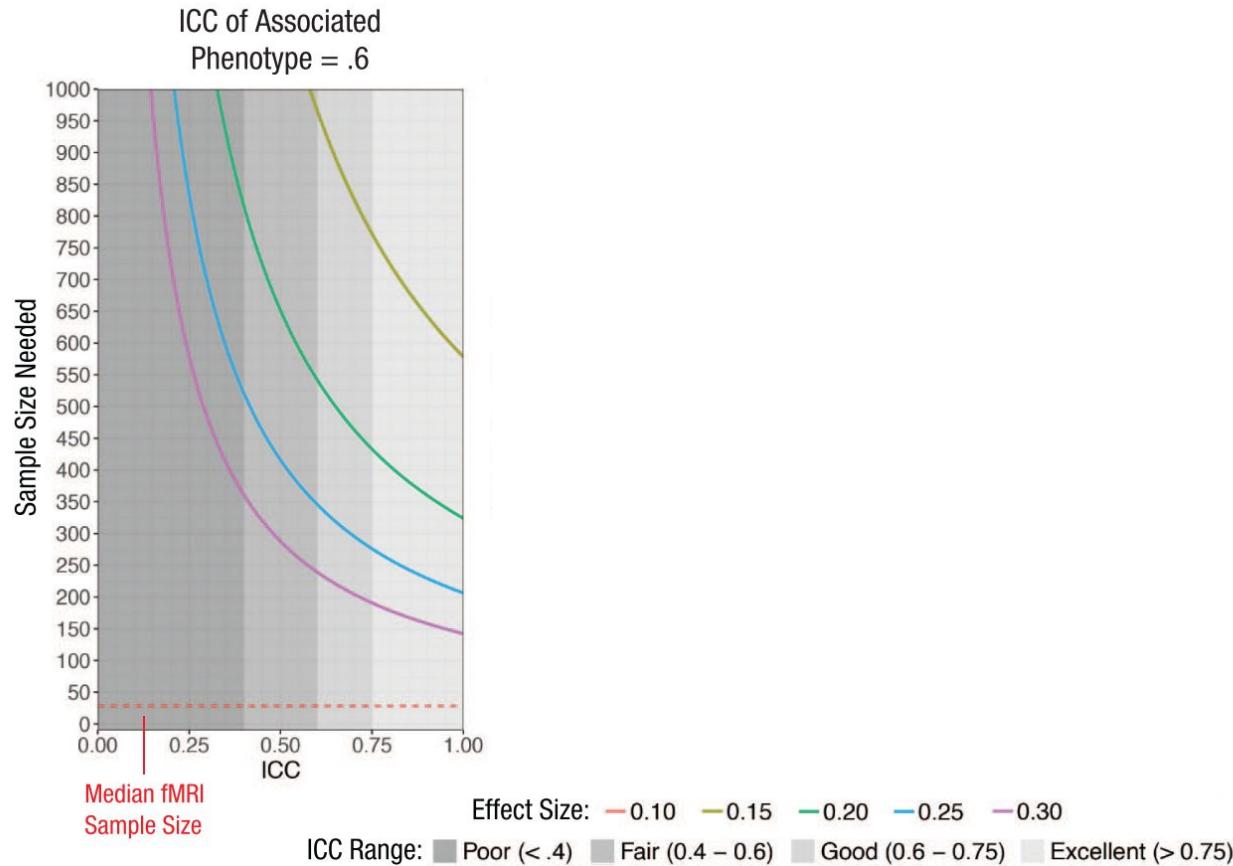
Maxwell L. Elliott<sup>1</sup>, Annchen R. Knodt<sup>1</sup>, David Ireland<sup>2</sup>,  
Meriwether L. Morris<sup>1</sup>, Richie Poulton<sup>2</sup>, Sandhya Ramrakha<sup>2</sup>,  
Maria L. Sison<sup>1</sup>, Terrie E. Moffitt<sup>1,3,4,5</sup>, Avshalom Caspi<sup>1,3,4,5</sup>,  
and Ahmad R. Hariri<sup>1</sup>



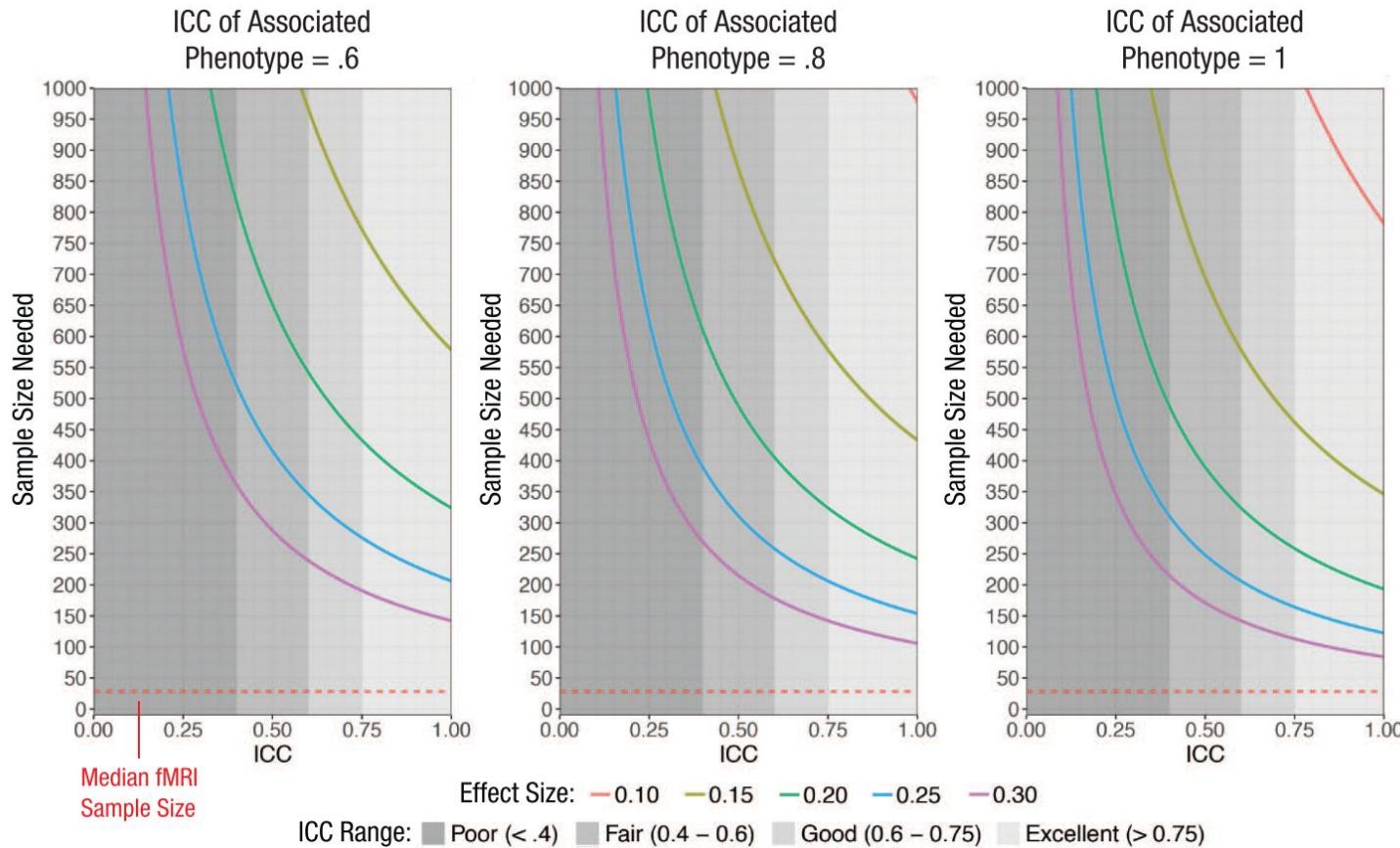
Psychological Science  
1–15  
© The Author(s) 2020  
Article reuse guidelines:  
[sagepub.com/journals-permissions](https://sagepub.com/journals-permissions)  
DOI: 10.1177/0956797620916786  
[www.psychologicalscience.org/PS](http://www.psychologicalscience.org/PS)  

**Fig. 1.** The influence of task-functional MRI (fMRI) test-retest reliability on the sample size required for 80% power to detect brain-behavior correlations of effect sizes commonly found in psychological research. Power curves are shown for three levels of reliability of the associated behavioral or clinical phenotype. The figure was generated using the *pwr.r.test* function in R (Champely, 2018), with the value for  $r$  specified according to the attenuation formula in the Appendix. ICC = intraclass correlation coefficient.



**Fig. 1.** The influence of task-functional MRI (fMRI) test-retest reliability on the sample size required for 80% power to detect brain-behavior correlations of effect sizes commonly found in psychological research. Power curves are shown for three levels of reliability of the associated behavioral or clinical phenotype. The figure was generated using the *pwr.r.test* function in R (Champely, 2018), with the value for  $r$  specified according to the attenuation formula in the Appendix. ICC = intraclass correlation coefficient.



**Fig. 1.** The influence of task-functional MRI (fMRI) test-retest reliability on the sample size required for 80% power to detect brain-behavior correlations of effect sizes commonly found in psychological research. Power curves are shown for three levels of reliability of the associated behavioral or clinical phenotype. The figure was generated using the *pwr.r.test* function in R (Champely, 2018), with the value for  $r$  specified according to the attenuation formula in the Appendix. ICC = intraclass correlation coefficient.

# Yet...

**fMRI can be highly reliable, but it depends on what you measure**

Philip A. Kragel<sup>1\*</sup>

Xiaochun Han<sup>2</sup>

Thomas E. Kraynak<sup>3</sup>

Peter J. Gianaros<sup>3</sup>

Tor D. Wager<sup>2\*</sup>

<sup>1</sup> Emory University

<sup>2</sup> Dartmouth College

<sup>3</sup> University of Pittsburgh

# Main ideas

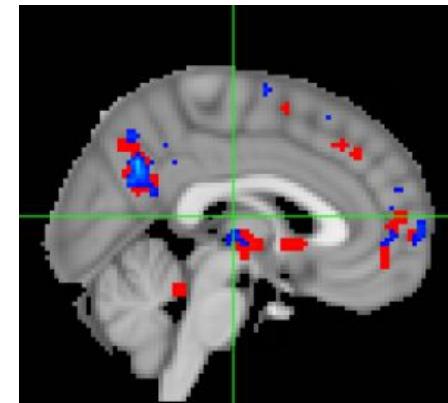
- Is machine learning just fancy statistics?
- Is machine learning better?
- How are people using ML in neuroscience?
- Things to consider when transitioning from voxelwise stats to ML

# It isn't quite this bad



# Often cannot interpret parameters as you're used to

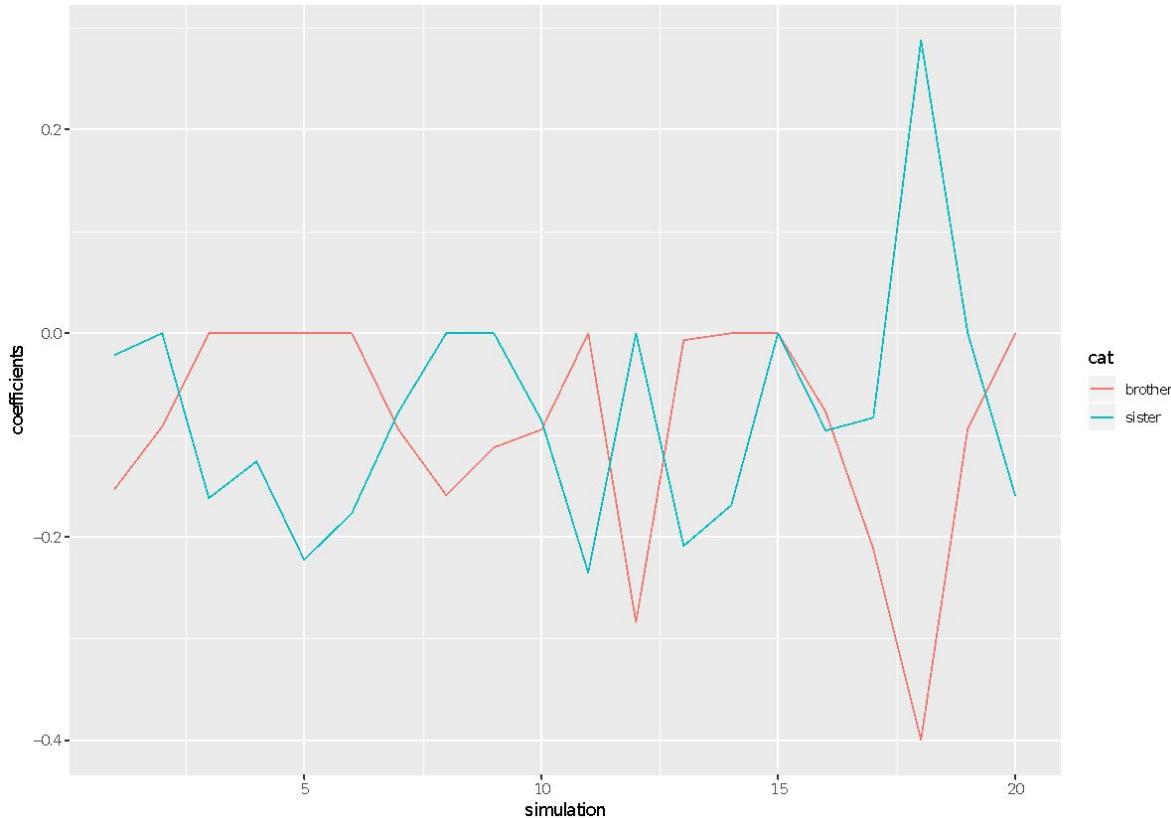
- Depends on how the regularization was done!
- The goal of the classifier is to classify
  - Doesn't necessarily matter what features it zeroed in on



# Example: Lasso regression

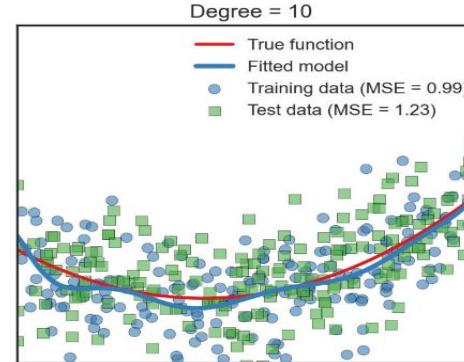
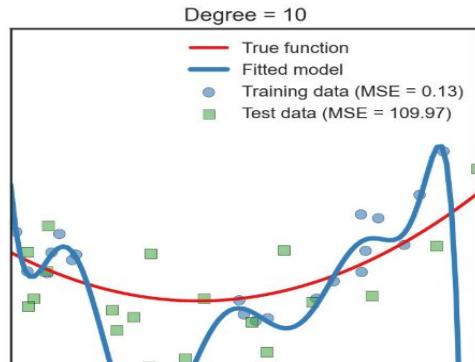
- Two cats are always right next to each other
- Their distance measures are correlated at 0.95
- Each cat's true regression parameter is -0.1
  - For every 10 feet of increase in distance, my happiness drops by 1 :(

# Lasso coefficients over multiple data sets



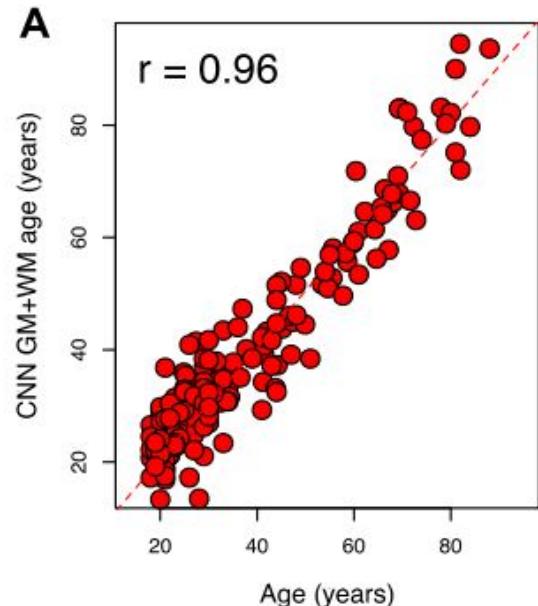
# You need *more* data

- “Big” data
  - Not wide (although tend to be wider)
  - Must have more observations (longer)

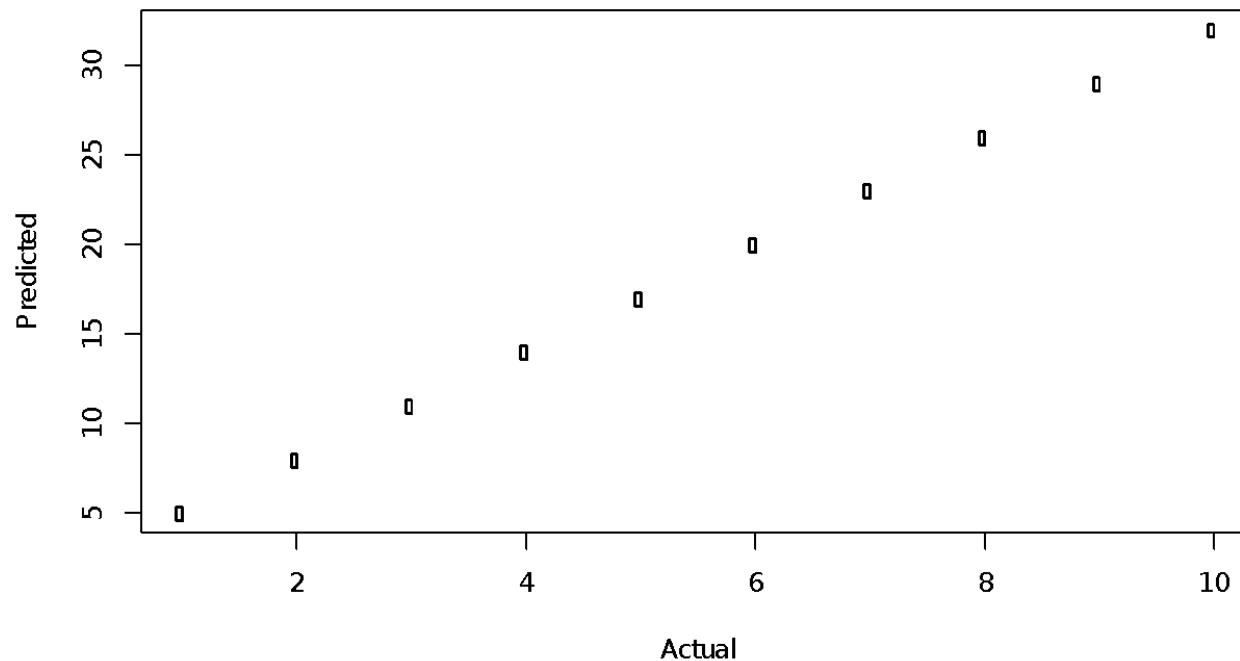


# Don't (only) use correlation to assess prediction of continuous measures

- Older folks' predictions are *older* and younger folks' are *younger*
- Correlation is high
  - They used plenty of other measures to assess



# Perfect correlation, bad prediction

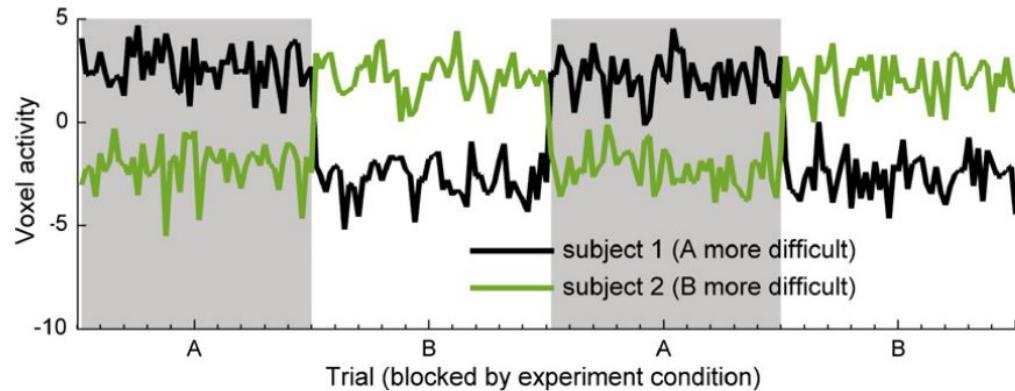


# Cross-validation isn't a cure-all

- Still can overfit data
- Need a left out validation set
  - You'll only validate *once*
- If you're developing new features, you may need a 3rd data set!

# Confounds impact analyses differently

Simulated example: experiment condition confounded with difficulty

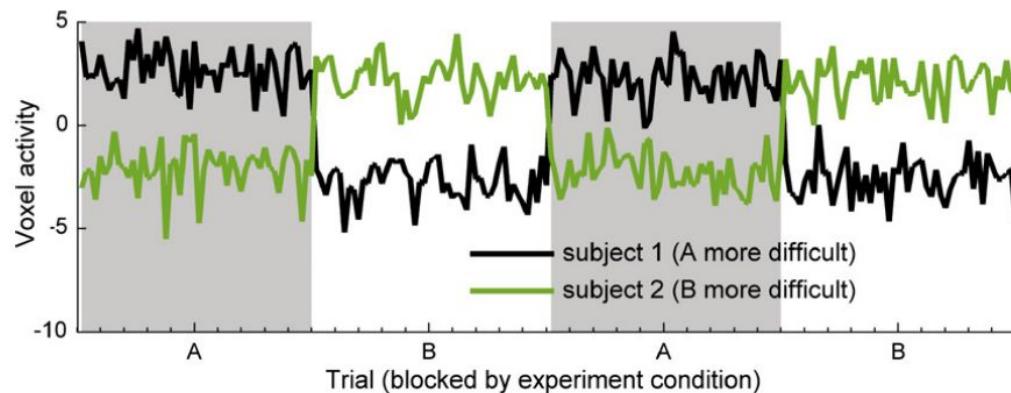


Example (I made this up): Compare math to reading, some people might find math more difficult

On average, difficulty effect should average out to reveal math/reading effect

# Confounds impact analyses differently

Simulated example: experiment condition confounded with difficulty



## Subject

Subject 1

Subject 2

## Individual-Subject Summary Statistics

### Experimental Effect (GLM)

$$\text{mean(A)} - \text{mean(B)} = +4.75$$

$$\text{mean(A)} - \text{mean(B)} = -5.56$$

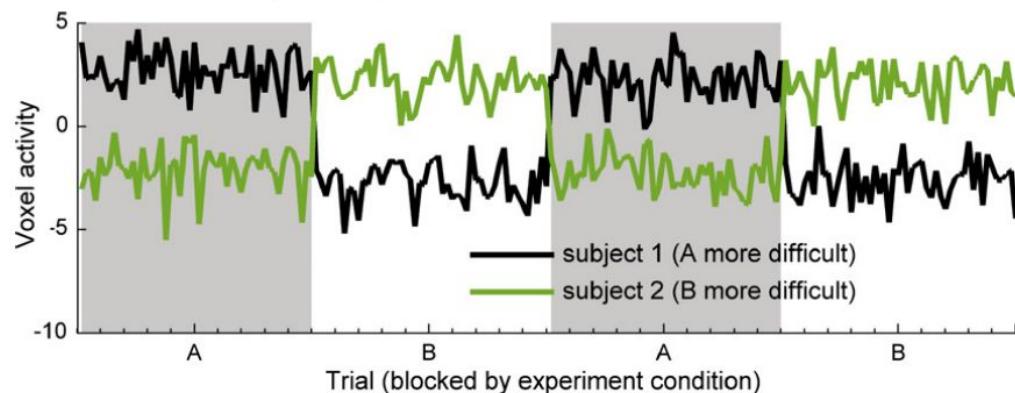
## Group Test Statistics (two-tailed t-test)

### Experimental Effect (GLM)

$$\begin{aligned} \text{mean(A)} - \text{mean(B)}: \\ t_1 = -0.0780, p = 0.9504, \text{n.s.} \end{aligned}$$

# Confounds impact analyses differently

Simulated example: experiment condition confounded with difficulty



Subject	Individual-Subject Summary Statistics	
	Experimental Effect (GLM)	Discrimination Success (MVPA)
Subject 1	$\text{mean(A)} - \text{mean(B)} = +4.75$	classification accuracy = +13.15, within-minus-across = +3.826
Subject 2	$\text{mean(A)} - \text{mean(B)} = -5.56$	classification accuracy = +13.44, within-minus-across = +3.848
Group Test Statistics (two-tailed <i>t</i> -test)		
Experimental Effect (GLM)	Discrimination Success (MVPA)	
	$\text{mean(A)} - \text{mean(B)}: t_i = -0.0780, p = 0.9504, \text{n.s.}$	classification accuracy: $t_i = 94.0, p < 0.01, \text{sig.}$
		within-minus-across: $t_i = 348, p < 0.01, \text{sig.}$

# Dive in! Ask questions and be careful



# Questions?